

"Express Mail" mailing label number:

EL071859778US

# METHOD AND SYSTEM FOR INDEXING INFORMATION AND PROVIDING RESULTS FOR A SEARCH INCLUDING OBJECTS HAVING PREDETERMINED ATTRIBUTES

Daniel Mattias Larsson  
Jan Olof Gelin  
Ian Douglas Hegerty

## BACKGROUND OF THE INVENTION

### Field of the Invention

The invention relates generally to data retrieval systems, and more specifically to a method and system for indexing information, searching information and providing search results responsive to a query.

### Description of the Related Art

The Internet has become one of the greatest information resources the world has ever known. Web sites throughout the world include a rich collection of text, images, audio, and video. However, finding specifically sought information among these web pages is not an easy task. Several search engines have been created to assist users in seeking information of interest. Most search engines feature searches of the text portion of web pages. These search engines typically include an index of web pages retrieved by its automated browsers (called robots or spiders). Thus, when a search engine receives a generated query including one or more keywords, the search engine locates pages that include at least one of the keywords and serves the user web pages that list the pages that include the one or more keywords received. The search results usually include a title and some information scraped from the web page. The user then selects the web pages to be displayed by clicking on the corresponding hyperlink. In many cases, many search results are servers and selection is not easy because the descriptive text shown in the search may not convey the full content of the corresponding page. There are also multimedia search services which locate documents containing or linking to multimedia files such as images in response to text search queries.

Examples of such search services include the AltaVista image search and the Lycos multimedia search for pictures. However, such services search only for documents containing images. Moreover those searches search only specialized image indexes. Therefore, there is a need for a search system wherein search results provide the user a better indication of the content of the web pages identified responsive to a received query.

## **SUMMARY OF THE INVENTION**

Briefly according to one embodiment of the invention, a system and method for providing at least one search result responsive to receiving a search query comprising at least one search query term, includes receiving the search query; locating at least one record containing at least one search query term; displaying a text representation of each of the records containing at least one search query term; and displaying a representation of an object in relation with the text representation for each record that includes an object comprising at least one predetermined attribute.

According to another embodiment of the invention, a method for indexing records in an index of an information network, comprises the acts of: receiving a record; searching the record for an indication that an object comprising at least one specified attribute is to be displayed with the record; finding the indication; creating a representation of the object; storing the representation in association with the record to which it corresponds; and making an entry for the record in the index, said entry indicating the object to be displayed with the record.

## **BRIEF DESCRIPTION OF THE DRAWINGS**

The present invention may be better understood, and its numerous objects, features and advantages made apparent to those skilled in the art by referencing the accompanying drawings. The use of the same reference number throughout the several figures designates a like or similar element.

FIG. 1 is a block diagram illustrating a distributed computer system for indexing a plurality of distributed information records according to an embodiment of the invention.

FIG. 2 is a block diagram of a search engine including an index.

FIG. 3 is a block diagram of pages parsed by the search engine of FIG. 2.

FIG. 4 is a search result page according to an aspect of the invention.

FIG. 5 is a flow chart illustrating a method according to an aspect of the invention.

FIG. 6 is a flow chart illustrating a method according to another aspect of the  
5 invention.

## **DETAILED DESCRIPTION**

FIG. 1 shows a distributed computer system 100 including a plurality of distributed  
information records 200 to be indexed. The distributed system 100 includes client computers  
110 connected to server computers (sites) 120 via a network 130. The network 130 can use  
10 Internet communications protocols (IP) to allow the clients 110 to communicate with the  
servers 120.

The client computers 110 can be personal computers, workstations, or larger or  
smaller computer systems. Each client computer 110 typically includes one or more  
processors, memories, and input/output devices. The servers 120 can be similarly configured.  
15 However, in many instances server sites 120 include many computers, perhaps connected by  
a separate private network. In fact, the network 130 may include hundreds of thousands of  
individual networks of computers. Although the client computers 110 are shown as being  
separate from the server computers 120, it should be understood that a single computer can  
perform both client and server roles.

20 During operation of the distributed system 100, users of the clients 110 desire to  
access information records 122 stored by the servers. The records of information 122 can be  
in the form of Web (or HTML) pages 200 or other data formats such as FTP (file transfer  
protocol). The pages 200 can contain or link to data records including as content plain textual  
information, or more complex digitally encoded multimedia content, such as software  
25 programs, images, graphics, audio signals, videos, and so forth.

It should be understood that although this description focuses on locating information  
on the World Wide Web (the web), the system can also be used for locating and indexing

information via other wide or local area networks (WANs and LANs), or information stored in a single computer using other communications protocols.

The client computers 110 can execute Web browser programs 112, such as NAVIGATOR, EXPLORER or MOSAIC (all trademarks of their respective owners) to download and display the pages or records 200. The browser programs 112 allow the users to enter addresses of specific Web pages 200 to be retrieved. Typically, the address of a Web page is specified as a Universal Resource Locator (URL). In addition, once a page has been retrieved, the browser programs 112 can provide access to other pages or records responsive to a user "clicking" on hyperlinks in previously retrieved Web pages. Such hyperlinks provide an automated way to enter the URL of another page, and to retrieve that page.

The search engine 140 assists users to identify pages of interest among the millions of pages which are available on the Web. The search engine 140 parses the pages 200, indexes the parsed pages, searches the index 70, and presents users with information about the pages 200 located.

The search engine 140 can be configured as one or more clusters of symmetric multi-processors (P) 142, for example, Digital Equipment Corporation ALPHA processors, memories (M) 144, disk storage devices 146, and network interfaces 148 are connected to each other by high speed communications buses 143. Although, the ALPHA processors 142 are 64 bit RISC processors, the search engine 140 can be any type of processor, which has sufficient processing power and memories, including 32 bit CISC processors. For smaller collections of information, the search engine can be run on the computer storing the database.

### Search Engine Overview

FIG. 2 shows the basic components of the search engine 140. The search engine 140 can include an automated Web browser or spider 20, a parsing module 30, an indexing module 40, a query module 50, index stream readers (ISR) 60, an index 70, and a maintenance module 80.

### Browsing

During the operation of the search engine 140, the spider 20 periodically sends out service requests 21 over the network 130. The requests 21 include URLs. In response to the

requests 21, the sites 120 return the records or pages 200 to the spider 20. The spider 20 can locate additional pages by following hyperlinks embedded in previously acquired pages. The browser 20 is described more completely in U.S. patent 5,974,455 entitled "System and Method for Locating Pages on the World Wide Web" and which is hereby incorporated by reference.

### **Parsing**

The pages 200 can be presented to the parsing module 30 as they are received or in batches which may amount to ten thousand pages or more, at one time. The parsing module 30 breaks down the portions of information of the pages 200 into fundamental indexable elements or atomic pairs 400. The word is a literal representation of the parsed portion of information, the location is a numeric value. The pages are parsed in order of the location of the words such that a location of the first word of a next page follows a location of the last word of a previous page. The parsing module 30 assigns increasing integer numbers to the locations, although other sequential orderings are also possible.

### **Indexing**

The indexing module 40 sorts the pairs 400, first in word order, and second in location order. The sorted pairs 400 are used to generate the index 70 of the words of the pages 200. Abstractly, the index 70 can be pictured as comprising compressed data structures 71 and summary data structures 72-73. The compressed data structure 71 is a compression of the word location pairs 400. The data structure 72 is a summary of the structure 71 and the data structure 73 is a summary of data structure 72. The structures 71 and 72 can be stored on disk, and the structure 73 can be stored in DRAM (dynamic random-access memory).

In the data structure 71, each word representing a unique portion of information of the pages 200 is stored only once. All of the locations which are instances of the word in the pages 200 are stored following the word. The locations follow the word in an order according to their locations. The locations essentially are pointers to the parsed portions of information.

The data structures of the index 70 are optimized for query access. This means that the word-location pairs 400 are compressed to reduce storage, and uncompressing is

minimized in order to preserve processor cycles during searching. Furthermore, the data structures of the index 70 also allow concurrent maintenance of the index 70 to delete old entries and to add new entries while queries are processed.

### **Querying**

5 Users interact with the index 70 via the query module 50 by providing queries 52. Users can be located remotely or locally with respect to the search engine 140. The terms of a query can include words and phrases, e.g., multiple words enclosed in quotation marks ("). The terms can be related by Boolean operators such as OR, AND, and NOT to form expressions.

10 During operation, the query module 50 analyzes the queries 52 to generate query requests 54. The query requests invoke a small number of basic types of object-oriented index stream readers (ISRs) 60. The index stream readers 60 sequentially scan the data structures 71-73 in a manner to minimize the amount of data that need to be uncompressed.

15 As a result of searching the index 70 by the stream reader objects 60, addresses 56 of pages which are qualified by the queries are identified. A presentation module 58 delivers information 59 about the qualifying pages to the users. The information 59 can include a summary of the pages located. Using the summary information, the users can access the identified pages with Web browsing software, or other techniques. The information 59 is preferably provided in the form of Web pages including a list of sites satisfying the received query. Each search result preferably includes a list of sites satisfying the received query. Each search result preferably includes a hyperlink to the subject site.

### **Maintaining**

25 The maintenance module 80 is used to add and delete information of the index 70. Modified pages can be handled as a delete and add operation. A particular problem solved is to allow substantially continuous access to the index 70 by millions of users each day as the index 70 is concurrently updated. The maintenance module 80 also effectively deals with duplicate Web pages containing substantially identical content.

## Indexing and Serving Images and Other Objects

According to one embodiment of the invention, the spider 20 searches for indications (e.g., links) that images 121 (or other non-markup objects such as audio files) having predetermined attributes or characteristics such as a high quality image (determined by size, aspect ratio, and number of colors) are to be displayed with the subject Web page. Once found the qualifying image is extracted, converted to a thumbnail (reduced resolution) version and indexed along with the parent page identifier (URL). Alternatively, the image may be stored in full-resolution form and converted to a thumbnail for serving to users as part of the search results. The indexing may be done simply by including an entry in the index 70 (e.g., a unique identifier) that indicates that the subject page includes a link to an image comprising the predetermined attributes. The qualifying images are stored to be transmitted for display to end-users. The images may be stored within the system 140 or may be stored in any other server connected to the web 130 for serving to end-users therefrom.

## The Parsing Module

### Words

As shown in FIG. 3, the records or pages 200 are parsed by the parsing module 30 in the order that pages are received from the browser 20. The parsing module 30, in a collating order of the sequential locations of the content, breaks the information of the pages 200 down into discrete indexable elements or individual "words" 300. In the index 70 each word is stored as a "literal" or character based value. It should be understood, that the terms "page" 200 and "word" 300 are used to represent many different possible content modalities and data record specifications. According to an embodiment, the parser finds any indication within the subject page that an image is to be displayed therewith.

### Pages

A page 200 can be defined as a data record including a collection of portions of information or "words" having a common database address, e.g., a URL. This means that a page can effectively be a data record of any size, from a single word, to many words, e.g., a large document, a data file, a book, a program, or a sequence of images.

On the Web, it has become common to represent information using a hypertext Markup Language (html). In this case, the pages can include other "marks" which indicate how the "words" of the page are to be processed and presented. Pages can include programs, for example JAVA applets, which may require specialized parsing. The information of some  
 5 pages can be expressed in a programming

language, for example, Postscript (.ps), or Acrobat (.pdf) files. The pages 200 can encode multimedia items including digitized images, graphic, audio or video components.

The pages or data records 200 are not necessarily Web pages. For example, the pages can be composed of portions of information of other databases, for example, all of the case  
 10 law in the United States. Even if such pages do contain hyperlinks, they may contain other types of links. In this context, the links mean references in one document which can be used to find other documents. Although hyperlinks are one example, many other types of links may be processed.

#### **Word and Location Pairs**

Each page 200 is broken down into a sequence of pairs 400 according to the collating  
 15 order of the locations of the words 300. Each pair 400 stores a word and its location. The locations of the words indicate the relative order in which the parsing module identified the words 300 in the pages 200.

#### **Attributes and Metawords**

20 In addition to recognizing locations and words, the parsing module 30 also detects and encodes attributes about the content of the records or pages. Attributes can be associated with entire pages, portions of pages 230, 240, 250, and 260, e.g., fields, or individual words 203.

||| Attribute values, as defined herein, are expressed as "metawords." Metawords are also stored as literals, this means that the search engine 140 treats metawords the same as words.

25 Therefore, a metaword is associated with a location to form a pair metaword, location. Thus, according to an embodiment of the invention, for pages wherein an image or other object is located, a metaword can be indexed indicating the presence of the image or object in the document.



For example, the page 200 of FIG. 3 can have several associated attributes. Portion 230 has an attribute associated with it that indicates the nature of the text in that portion. Portion 240 can be selections from various audio files. Portion 250 includes an attribute indicating selections of video files and portion 260 includes an attribute indicating an image.

5

### Serving Blended Search Results

Referring to FIG. 4, there is shown a search results page served in response to the query "Tony Blair." These results are "blended" in that some results include only text abstracts and others include both the abstract and a thumbnail image associated therewith.

10 Results 1, 2, and 3 only include a title, a short excerpt from the text of the subject page, the page URL, the date last modified, the language of the page, and user options (e.g., to save to a "my search" file or to email the page to a friend). That is because those results do not include images including predetermined attributes. In addition to the above items result 4 shows an image to which the subject page includes a link. The image is shown such that the  
15 end user can associate the image with the relevant result. In this case the image is shown alongside the text representation of the content of the subject page but the image may be shown underneath the text representation or elsewhere in the result page as long as some indication is made of the search result to which it relates.

20 The image shown is preferably the first image appearing in the page that meet the predetermined criteria (e.g., it comprises certain attributes or characteristics). However, it is also possible to display the image most relevant to the query by scanning annotations or other text in the subject page relating to the image.

25 In one embodiment, only the first image appearing in a crawled page is retrieved, stored, and indexed. However, it is also possible to retrieve and store every image found in every page, and index therewith certain metawords relevant to the query may be served.

Referring to FIG. 5, there is illustrated a method 500 for providing at least one search result responsive to receiving a search query including at least one search term (or keyword). The method 500 includes the acts of: receiving the search query (502); locating at least one qualifying record (e.g., a Web page) containing at least one search query term (504);

transmitting a text representation (e.g., an abstract) of each of the qualifying records for display to a user (506); and transmitting a representation of an object (e.g., a thumbnail image) in relation with the text representation of each qualifying record that includes an object comprising at least one predetermined attribute for display to the user (508). Search results that do not include images are presented as text-only descriptions of the pages to which they relate.

In the preferred embodiment, each thumbnail image is displayed next to the text representation of the record to which it corresponds. In other embodiments, the search engine can index references to other objects such as programs, audio files, video files or other multimedia files and store representations of such objects for serving to users as a part of search results.

#### **Method for Indexing Images to be Displayed**

Referring to FIG. 6, there is illustrated a method 600 for indexing records in an index of an information network. The method 600 comprises the acts of: receiving a page or other record (602); determining whether the page includes a link to an image, or other object, comprising at least one specified attribute such as a predetermined size (604); determination 604 is affirmative, the image is indexed (608) by making an entry in the search engine index that indicates the presence of the link to the image in the subject page; and a representation of the image such as a thumbnail (a reduced resolution version) is created and stored (610) for serving to users. Thus, according to this embodiment of the invention the spider 20 is modified to mine the retrieved documents for images to be indexed so that they can be served with search results to which they pertain.

In the preferred embodiment images having predetermined attributes are indexed and stored to be served with search results. However, as in the case of the search procedure 500, other objects may be located within Web pages crawled and those objects can be indexed and stored for serving to end users.